

(1) (a) since y is a one-hot vector with
1 at word 0 and 0 elsewhere,

$$- \left[y^T \log(\hat{y}) \right] = - \left[\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) \right]$$

$$= - \left[\sum_{j \neq 0} (0) (\log(\hat{y}_j)) + (1) \log(\hat{y}_0) \right]$$

$$= - \log(\hat{y}_0)$$

$$(1) (b) \mathcal{J}_{\text{naive softmax}}(v_c, 0, V) = -\log \left[\frac{e^{u_0^T v_c}}{\sum_w e^{u_w^T v_c}} \right]$$

$$\frac{\partial \mathcal{J}}{\partial v_c} = \frac{\partial \mathcal{J}}{\partial \hat{y}_0} \cdot \frac{\partial \hat{y}_0}{\partial v_c} = -\frac{1}{\hat{y}_0} \left[\frac{\sum_w e^{u_w^T v_c} (u_0 e^{u_0^T v_c})}{\left(\sum_v e^{u_v^T v_c} \right)^2} \right]$$

$$- \frac{e^{u_0^T v_c} \left(\sum_k v_c e^{u_k^T v_c} \right)}{\left(\sum_v e^{u_v^T v_c} \right)^2}$$

$$= -\frac{1}{\hat{y}_0} \left[\frac{\sum e^{u_w^T v_c} (u_0 e^{u_0^T v_c}) - e^{u_0^T v_c} \left(\sum_k u_k e^{u_k^T v_c} \right)}{\left(\sum_v e^{u_v^T v_c} \right)^2} \right]$$

$$= -\frac{1}{\hat{y}_0} \left[(1)(u_0 \hat{y}_0) - (\hat{y}_0) \left(\sum_k u_k \hat{y}_k \right) \right]$$

$$= -u_0 + \sum_k u_k \hat{y}_k = -\sum_i y_i u_i + \sum_k u_k \hat{y}_k$$

$$= -V^T y + V^T \hat{y} = V^T (\hat{y} - y)$$

$$(1) (c) \quad \frac{\partial J}{\partial u_w} = \frac{\partial J}{\partial \hat{y}_w} \cdot \frac{\partial \hat{y}_w}{\partial u_w}$$

$$\text{Let } w=0, \quad \left(-\frac{1}{\hat{y}_0} \right) \left[\frac{\left(\sum_w e^{u_w^T v_c} \right) (v_c e^{u_0^T v_c}) - (e^{u_0^T v_c}) (v_c \sum_w e^{u_w^T v_c})}{\left(\sum_w e^{u_w^T v_c} \right)^2} \right]$$

$$= -\frac{1}{\hat{y}_0} \left[\frac{\left(\sum_w e^{u_w^T v_c} \right) (v_c e^{u_0^T v_c}) - (e^{u_0^T v_c}) (v_c \sum_w e^{u_w^T v_c})}{\left(\sum_w e^{u_w^T v_c} \right)^2} \right]$$

$$= -\frac{1}{\hat{y}_0} \left[(1)(v_c)(\hat{y}_0) - v_c (\hat{y}_0)^2 \right] = v_c (\hat{y}_0 - 1)$$

$$\text{Let } w \neq 0, \quad w=z, \quad -\frac{1}{\hat{y}_0} \left[\frac{\left(\sum_w e^{u_w^T v_c} \right) (0) - (e^{u_z^T v_c}) (v_c e^{u_0^T v_c})}{\left(\sum_w e^{u_w^T v_c} \right)^2} \right]$$

$$= -\frac{1}{\hat{y}_0} \left[-\hat{y}_z v_c \hat{y}_0 \right] = \hat{y}_z v_c = v_c (\hat{y}_z - \frac{y}{z})$$

$$\therefore \frac{\partial J}{\partial u} = \underbrace{(\hat{y} - y)}_{N \times 1} \underbrace{v_c^T}_{1 \times D}$$

$$(1) (d) \sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$$

$$\sigma'(x) = \frac{(e^x+1)(e^x) - (e^x)(e^x)}{(e^x+1)^2}$$

$$= \sigma(x) - \sigma(x)^2 = \sigma(x)(1 - \sigma(x))$$

$$(1) (e) \quad J_{\text{neg-sample}}(v_c, o, U) = -\log(\sigma(u_o^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c))$$

$$\begin{aligned} \frac{\partial J}{\partial v_c} &= \frac{-1}{\cancel{\sigma(u_o^T v_c)}} \cdot \left[\cancel{\sigma(u_o^T v_c)} \cdot [1 - \sigma(u_o^T v_c)] \right] (u_o) \\ &\quad - \sum_{k=1}^K \frac{1}{\cancel{\sigma(-u_k^T v_c)}} \cdot \cancel{\sigma(-u_k^T v_c)} \cdot [1 - \sigma(-u_k^T v_c)] (-u_k) \\ &= -u_o (1 - \sigma(u_o^T v_c)) + \sum_{k=1}^K u_k (1 - \sigma(-u_k^T v_c)) \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial u_o} &= \frac{-1}{\sigma(u_o^T v_c)} \cdot \sigma(u_o^T v_c) \cdot (1 - \sigma(u_o^T v_c)) v_c \\ &= -v_c [1 - \sigma(u_o^T v_c)] \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial u_k} &= \frac{-1}{\cancel{\sigma(-u_k^T v_c)}} \cdot \cancel{\sigma(-u_k^T v_c)} \cdot (1 - \sigma(-u_k^T v_c)) (-v_c) \\ &= v_c [1 - \sigma(-u_k^T v_c)] \end{aligned}$$

use for-loop for neg samples. Possible Duplicates

Negative sampling is more efficient than naive softmax. When computing $\frac{\partial J}{\partial v_c}$, naive softmax uses the entire word matrix U , whereas negative sampling uses only the target outside word vector, u_o and K negative samples u_1, \dots, u_K to compute gradients.

$$(1) (f) \mathcal{J}_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \mathcal{J}(v_c, w_{t+j}, U)$$

$$(i) \frac{\partial \mathcal{J}_{SG}}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial \mathcal{J}(v_c, w_{t+j}, U)}{\partial U}$$

$$(ii) \frac{\partial \mathcal{J}_{SG}}{\partial v_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial \mathcal{J}(v_c, w_{t+j}, U)}{\partial v_c}$$

$$(iii) \frac{\partial \mathcal{J}_{SG}}{\partial v_w} = 0$$

$w \neq c$